# Minería de datos y Patrones

Version 2024-1

## Selección de Características

[ Capítulo 3 ]

## Dr. José Ramón Iglesias

DSP-ASIC BUILDER GROUP

Director Semillero TRIAC

Ingenieria Electronica

Universidad Popular del Cesar

# 5 reasons why to select features

1. To avoid non-discriminative features
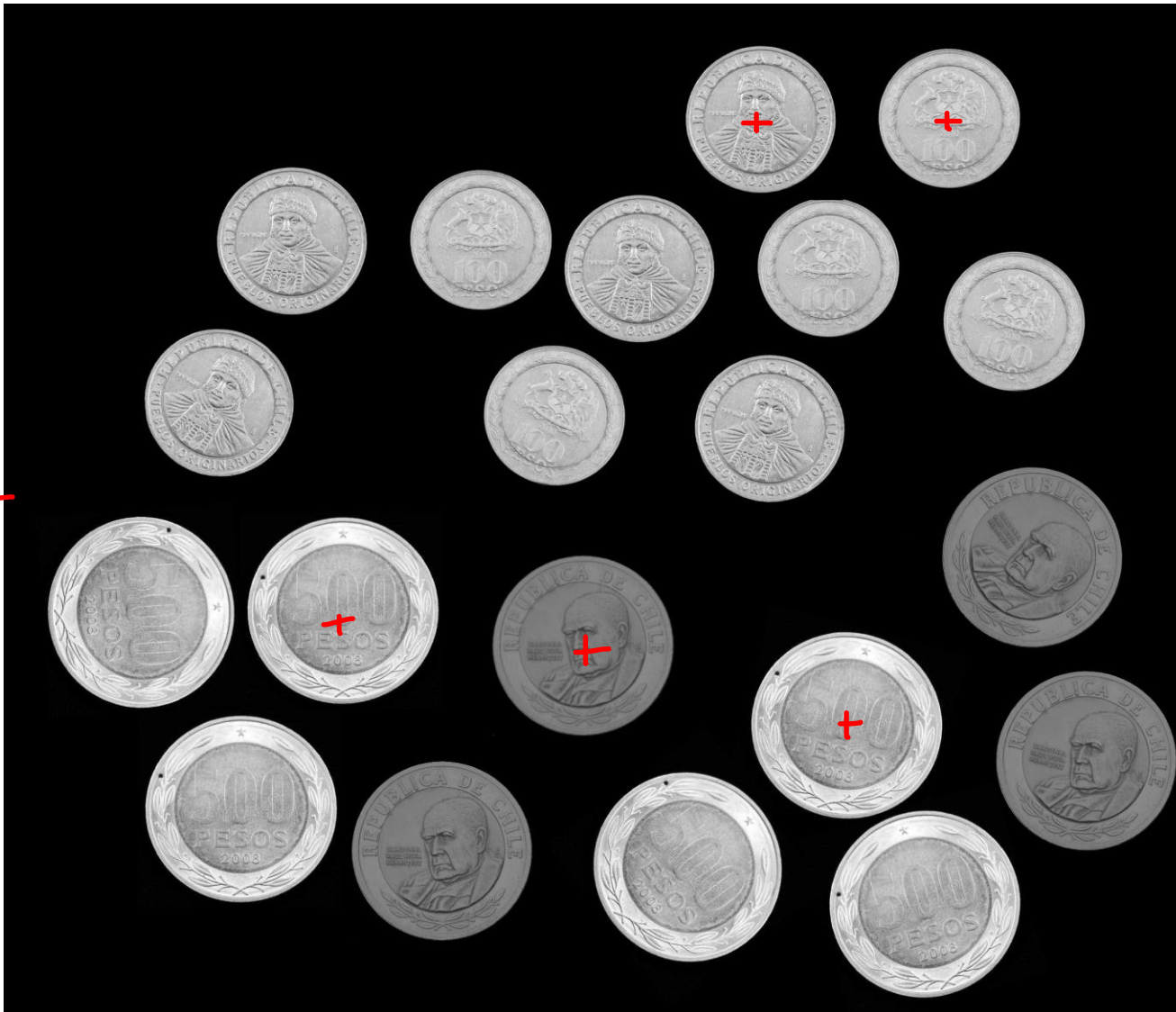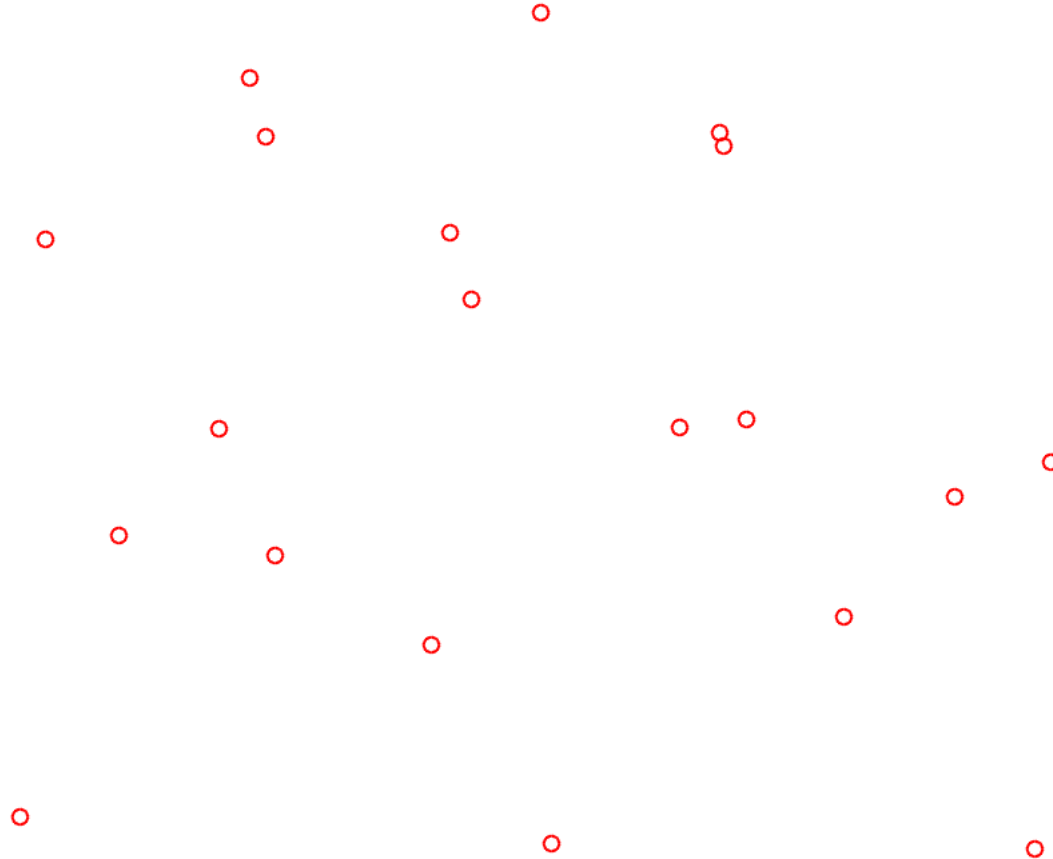2. To avoid correlated features
3. To simplify the testing stage
4. To avoid false correlations
5. To avoid the curse of dimensionality

# Example: Classification of Coins

# 5 reasons why to select features

1. To avoid non-discriminative features

# 1. To avoid non-discriminative features – ~~Grayvalue average~~ ☹

# 5 reasons why to select features

1. To avoid non-discriminative features
2. To avoid correlated features

## 2. To avoid correlated features – Area and ~~Diameter~~ ☹

# 5 reasons why to select features

1. To avoid non-discriminative features
2. To avoid correlated features
3. To simplify the testing stage

# 3. To simplify the testing stage – Area vs. ~~Hu+Flusser~~

# 5 reasons why to select features

1. To avoid non-discriminative features
2. To avoid correlated features
3. To simplify the testing stage
4. To avoid false correlations

# 4. To avoid false correlations – Location ☹

# 5 reasons why to select features

1. To avoid non-discriminative features
2. To avoid correlated features
3. To simplify the testing stage
4. To avoid false correlations
5. To avoid the curse of dimensionality

# Curse of dimensionality – 1D

(20 points in 1D)

# Curse of dimensionality – 2D



(20 points in 2D)

# Curse of dimensionality – 3D

When the dimensionality increases, the volume of the space increases so fast that the available data become sparse. This sparsity is problematic for any method that requires statistical significance.

[Wikipedia]

(20 points in 3D)
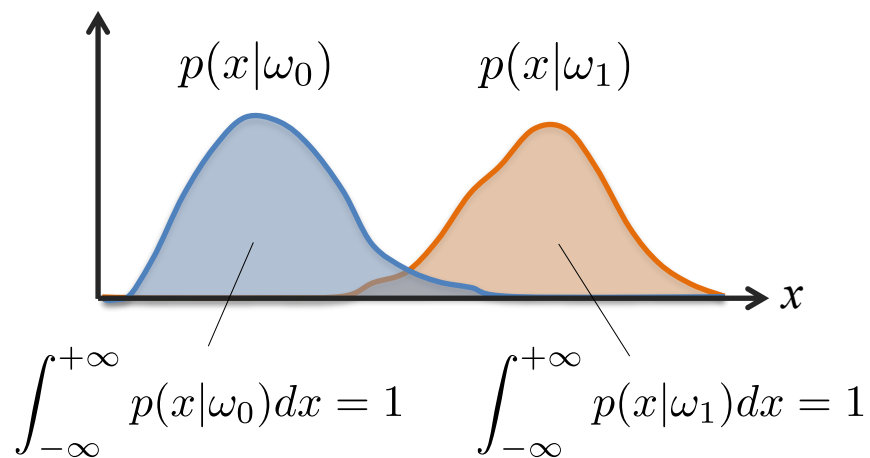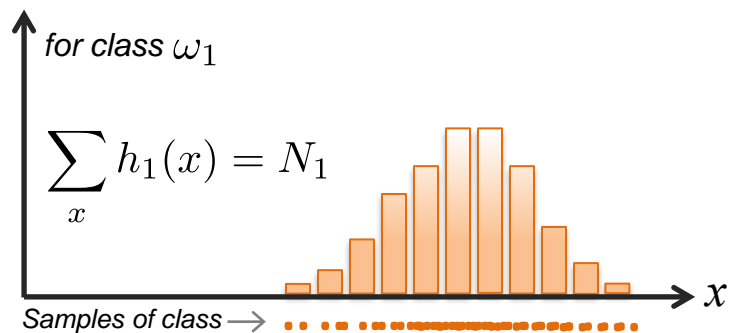
*Extracted features*

*Selected features*

*m features*

*labels*

*p features*

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_5$   $x_6$   $x_7$   $x_8$   $x_9$   $x_{10}$

$d$

$z_1$   $z_2$   $z_3$

*N samples*

*N samples*

**X**

**d**

**Z**

*Frequency distributions*

*for class $\omega_0$*

$$\sum_x h_0(x) = N_0$$

*for class $\omega_1$*

$$\sum_x h_1(x) = N_1$$

*Probability density functions*

$p(x|\omega_0)$    $p(x|\omega_1)$

$$\int_{-\infty}^{+\infty} p(x|\omega_0)dx = 1 \qquad \int_{-\infty}^{+\infty} p(x|\omega_1)dx = 1$$

*Frequency distributions*

*for class $\omega_0$*

$$\sum_x h_0(x) = N_0$$

$\leftarrow$ *Samples of class*

*for class $\omega_1$*

$$\sum_x h_1(x) = N_1$$

*Samples of class* $\rightarrow$

*Probability density functions*

$p(x|\omega_0)$

$p(x|\omega_1)$

$$\int_{-\infty}^{+\infty} p(x|\omega_0)dx = 1$$

$$\int_{-\infty}^{+\infty} p(x|\omega_1)dx = 1$$

Bad separability $x_1$

Good separability $x_2$

Very good separability $x_3$

score = LOW

score = MEDIUM

score = HIGH

*Bad*

$x_2$

$x_1$

*Good*

$x_2$

$x_1$

*Very good*

$x_2$

$x_1$

*Bad*  *Good*  *Very good*
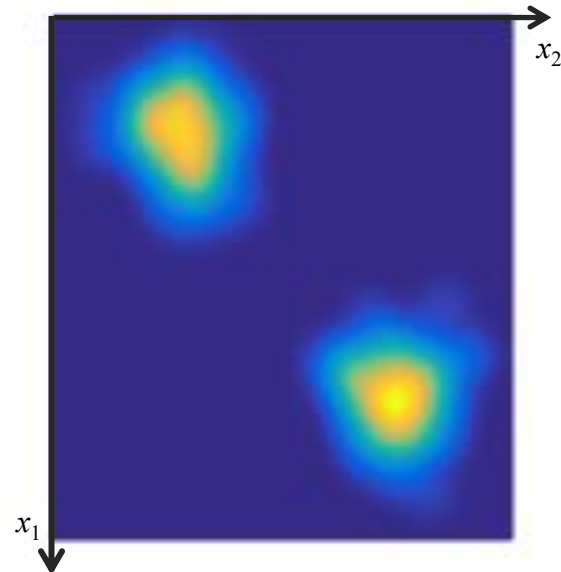
score = LOW

score = MEDIUM

score = HIGH

$\omega_0$
$\omega_1$
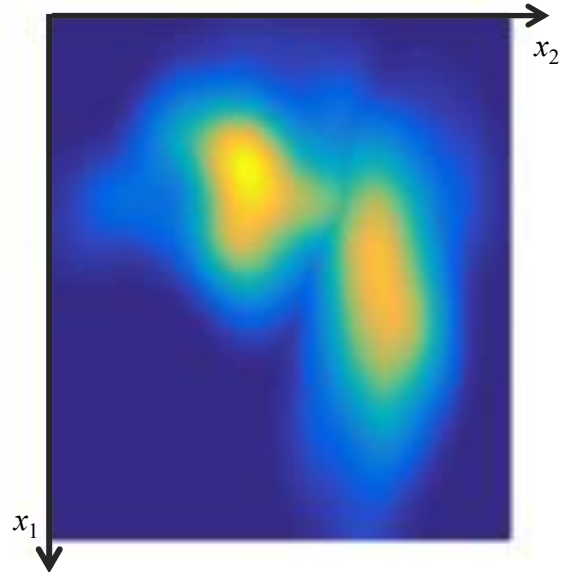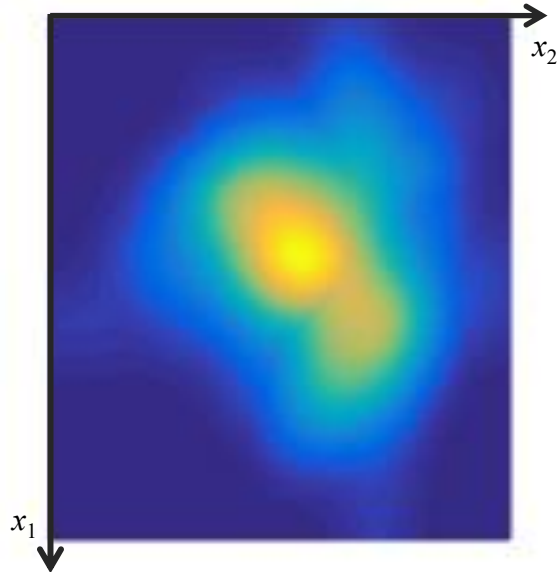
*Selected features <u>must be</u>…*

Can our objects of interest be everywhere?  —*yes*→  *invariant to translation*

Can our objects of interest be in any orientation?  —*yes*→  *invariant to rotation*

Can our objects of interest be of any size?  —*yes*→  *invariant to scale*